

P.Ż

---

# WPROWADZENIE DO ODCZYTYWANIA I ANALIZY DANYCH

**CO DZISIAJ – POBIERANIE TREŚCI ABY JĄ  
PÓŹNIEJ ANALIZOWAĆ – PODSTAWY PRZED  
BUDOWĄ WIĘKSZEJ SIECI. ODCZYTYWANIE  
TREŚCI Z RÓŻNYCH ŹRÓDEŁ**

# A MOŻE INACZEJ – ANALIZA TEKSTU

- ▶ Python jest językiem wykorzystywanym często we wszelkiej analizie danych
- ▶ Parsowanie tekstu w Pythonie to proces analizowania tekstu, którego celem jest zrozumienie jego struktury i znaczenia. Technika ta jest powszechnie stosowana w działaniach na danych, np. przy analizie logów, parsowaniu stron internetowych itp.
- ▶ W Pythonie istnieją różne biblioteki i moduły do parsowania tekstu, takie jak na przykład `re` (Regular Expressions) lub `BeautifulSoup`. Moduł `re` umożliwia parsowanie tekstu przy użyciu wyrażeń regularnych, zajmuje się on dopasowywaniem określonych wzorców w tekście. Natomiast `BeautifulSoup` służy do parsowania HTML i XML.
- ▶ Ogólnie rzecz biorąc, parsowanie tekstu w Pythonie polega na wyodrębnieniu określonych elementów tekstu na podstawie określonych kryteriów, np. wzorców wyrażeń regularnych. Znajomość podstawowych technik parsowania tekstu jest bardzo cenna i przydatna w różnych dziedzinach programowania.

# BEAUTIFUL SOUP TRANSFORMUJE DOKUMENTY HTML DO POSTACI DRZEWA OBIEKTÓW PYTHON

Cyt za dokumentacją modułu

# A MOŻE INACZEJ – ANALIZA TEKSTU

- ▶ BS4 wykorzystuje się klasycznie w połączeniu z request
- ▶ W trakcie naszej pracy, będziemy wykorzystywać bibliotekę BeautifulSoup do przetwarzania strony internetowej. Do tego celu wykorzystamy obiekt klasy BeautifulSoup, który będzie reprezentował stronę. Jako argument do konstruktora podamy tekst html strony, który wcześniej pobierzemy przy użyciu funkcji get z biblioteki requests. W drugim argumencie konstruktora podamy rodzaj parsera. BeautifulSoup to narzędzie nie tylko do parsowania HTML ale też XML. Na końcu ustalimy kodowanie w jednej z linii kodu. Przed przystąpieniem do przetwarzania strony, na sąsiedniej karcie, możemy zobaczyć kod HTML, który będziemy analizować."

```
● 1 import requests as req
   2 from bs4 import BeautifulSoup
```

```
1 import requests as req
2 from bs4 import BeautifulSoup
3 strona=req.get('http://jsystems.pl/static/data/pnl/dane.html')
4 strona.encoding='utf-8'
5 zupa=BeautifulSoup(strona.text,'html.parser')
6
```

# A MOŻE INACZEJ – ANALIZA TEKSTU Z WWW

- ▶ Metoda ".h3" wyodrębnia z obiektu "zupa" pierwszy znaleziony element HTML o nazwie "h3".
- ▶ - Funkcja "print" wyświetla ten znaleziony element HTML, który jest obiektem klasy "Tag" w przypadku znalezienia takiego elementu na stronie.
- ▶ - Funkcja "type" zwraca typ obiektu "zupa.h3", który w tym przypadku powinien być również "Tag".

```
1 import requests as req
2 from bs4 import BeautifulSoup
3 strona=req.get('http://jsystems.pl/static/data/pnl/dane.html')
4 strona.encoding='utf-8'
5 zupa=BeautifulSoup(strona.text,'html.parser')
6 print(zupa.h3)
7 print(type(zupa.h3))
```

```
1 import requests as req
2 from bs4 import BeautifulSoup
3 strona=req.get('http://jsystems.pl/static/data/pnl/dane.html')
4 strona.encoding='utf-8'
5 zupa=BeautifulSoup(strona.text,'html.parser')
6 print(zupa.find('h3'))
7 print(type(zupa.find('h3')))
```



# A MOŻE INACZEJ – ANALIZA TEKSTU Z WWW

- ▶ Im głębiej w strukturę strony wchodzimy, tym możemy więcej elementów używać

```
1 import requests as req
2 from bs4 import BeautifulSoup
3 strona=req.get('http://jsystems.pl/static/data/pnl/dane.html')
4 strona.encoding='utf-8'
5 zupa=BeautifulSoup(strona.text,'html.parser')
6 print(zupa.find('h3'))
7 print(type(zupa.find('h3')))
```

```
1 import requests as req
2 from bs4 import BeautifulSoup
3 strona=req.get('http://jsystems.pl/static/data/pnl/dane.html')
4 strona.encoding='utf-8'
5 zupa=BeautifulSoup(strona.text,'html.parser')
6 print(zupa.find(id='sekcja2'))
```

```
<div class="podsekcja" id="sekcja2">
    Zawartość sekcji numer 2
</div>
```

# A MOŻE INACZEJ – ANALIZA TEKSTU Z WWW

- ▶ Możliwości wyszukiwania jest więcej

```
1 import requests as req
2 from bs4 import BeautifulSoup
3 strona=req.get('http://jsystems.pl/static/data/pnl/dane.html')
4 strona.encoding='utf-8'
5 zupa=BeautifulSoup(strona.text,'html.parser')
6 print(zupa.find(class_='podsekcja'))
```

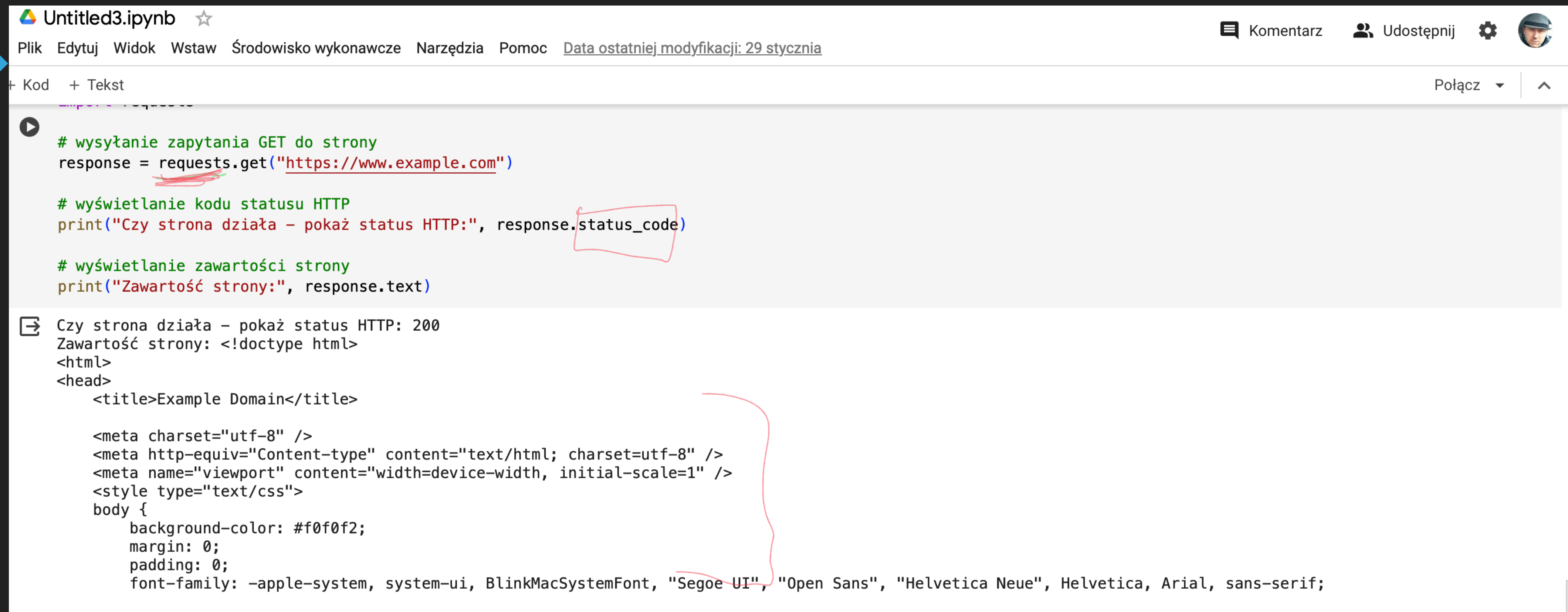
```
1 import requests as req
2 from bs4 import BeautifulSoup
3 strona=req.get('http://jsystems.pl/static/data/pnl/dane.html')
4 strona.encoding='utf-8'
5 zupa=BeautifulSoup(strona.text,'html.parser')
6 print(zupa.title)
7 print(zupa.head)
8 print(zupa.body)
```

```
1 import requests as req
2 from bs4 import BeautifulSoup
3 strona=req.get('http://jsystems.pl/static/data/pnl/dane.html')
4 strona.encoding='utf-8'
5 zupa=BeautifulSoup(strona.text,'html.parser')
6 print(zupa.find(attrs={'name':'stopka'}))
```



## A MOŻE INACZEJ – ANALIZA TEKSTU Z WWW

- ▶ Możliwości wyszukiwania jest więcej



The screenshot shows a Jupyter Notebook titled 'Untitled3.ipynb'. The interface includes a top menu bar with options like 'Plik', 'Edytuj', 'Widok', 'Wstaw', 'Środowisko wykonawcze', 'Narzędzia', 'Pomoc', and a date 'Data ostatniej modyfikacji: 29 stycznia'. On the right, there are icons for 'Komentarz', 'Udostępnij', and a user profile. Below the menu, there are tabs for 'Kod' and 'Tekst', and a 'Połącz' button. The main area contains a Python script with three comments: '# wysyłanie zapytania GET do strony', '# wyświetlanie kodu statusu HTTP', and '# wyświetlanie zawartości strony'. The script uses the 'requests' library to get a response from 'https://www.example.com' and prints the status code and text. The output shows the status code 200 and the HTML content of the page, which includes a title 'Example Domain' and various meta tags and CSS styles. Red annotations highlight the URL in the script, the 'status\_code' attribute in the print statement, and the 'Segoe UI' font name in the CSS output.

```
Untitled3.ipynb ☆
Plik Edytuj Widok Wstaw Środowisko wykonawcze Narzędzia Pomoc Data ostatniej modyfikacji: 29 stycznia
+ Kod + Tekst Połącz ^

▶ # wysyłanie zapytania GET do strony
response = requests.get("https://www.example.com")

# wyświetlanie kodu statusu HTTP
print("Czy strona działa – pokaż status HTTP:", response.status_code)

# wyświetlanie zawartości strony
print("Zawartość strony:", response.text)

⇒ Czy strona działa – pokaż status HTTP: 200
Zawartość strony: <!doctype html>
<html>
<head>
  <title>Example Domain</title>

  <meta charset="utf-8" />
  <meta http-equiv="Content-type" content="text/html; charset=utf-8" />
  <meta name="viewport" content="width=device-width, initial-scale=1" />
  <style type="text/css">
  body {
    background-color: #f0f0f2;
    margin: 0;
    padding: 0;
    font-family: -apple-system, system-ui, BlinkMacSystemFont, "Segoe UI", "Open Sans", "Helvetica Neue", Helvetica, Arial, sans-serif;
```

# ALTERNATYWY DLA BEAUTIFUL SOUP

---

- ▶ 1. lxml - biblioteka do przetwarzania i parsowania XML i HTML, która oferuje opcje szybkiego przetwarzania oraz układania DOM. Posiada również wiele funkcjonalności do przeszukiwania, transformowania i walidowania dokumentów XML i HTML.
- ▶ 2. PyQuery - biblioteka oparta na jQuery, która umożliwia przeszukiwanie i manipulowanie dokumentami HTML i XML. Jest to bardzo łatwa w użyciu biblioteka, która pozwala na korzystanie z wygodnego API do wyszukiwania elementów HTML i wykorzystywania selektorów CSS.
- ▶ 3. html5lib - biblioteka do przetwarzania dokumentów HTML5, która oferuje stabilne i dokładne parsowanie dokumentów. Jest ona nieco wolniejsza niż lxml, ale oferuje lepszą zgodność z HTML5 i może być wygodnym narzędziem do parsowania nieprawidłowo sformowanych dokumentów HTML.
- ▶ 4. re - biblioteka do pracy z wyrażeniami regularnymi, która może być używana do przetwarzania i parsowania tekstu HTML i XML. Jest to bardziej podstawowe podejście niż wykorzystanie biblioteki do parsowania dokumentów, ale może być przydatne w przypadku prostych zadań.
- ▶ 5. Scrapy - framework stanowiący kompletny zestaw narzędzi do przetwarzania i parsowania stron internetowych. Oferuje on elastyczne podejście do przetwarzania stron oraz bogatą funkcjonalność do przeszukiwania i pobierania danych z witryn internetowych, a także tworzenia pełnoprawnych aplikacji internetowych.

**WYKORZYSTUJEMY NP. PRZY  
PODEJMOWANIU DECYZJI ZWIĄZANYCH Z  
PARKINGAMI, AUTOMATYKĄ DŹWIĘKOWĄ ITP**



CZYM W OGÓLE JEST OCR

---

OCR TO ROZPOZNAWANIE ZNAKÓW LUB TEKSTU Z OBRAZÓW LUB RĘCZNIE STWORZONYCH DOKUMENTÓW.



**OCR TO ROZPOZNAWANIE ZNAKÓW LUB TEKSTU Z OBRAZÓW LUB RĘCZNIE STWORZONYCH DOKUMENTÓW.**

Do najpopularniejszych bibliotek w python, związanych z OCR należą:

- ▶ Tesseract - py te serve ct
- ▶ OCRopus
- ▶ Kraken
- ▶ easyOCR



## TESSERACT

Biblioteka, która w zasadzie jest API dostarczonym przez Google. Ogromne możliwości, ale raczej w większych zastosowaniach

```
1 import pytesseract (✓)
2 from PIL import Image pip install
3
4 # Wczytanie obrazu
5 image = Image.open('tekst.jpg')
6 # Konwersja obrazu na skalę szarości
7 image = image.convert('L') ← ?
8
9 # Rozpoznawanie tekstu przy użyciu Tesseract
10 text = pytesseract.image_to_string(image)
11
12 # Wypisanie rozpoznanego tekstu
13 print(text)
```

# KRAKEN

Od podstaw zbudowana pod kątem wykorzystywania sieci neuronowych (przykład rozbudowany później - zawiera więcej na początku niż Tesseract - 0 i 1 oznaczają, że ma być sam tekst )

```
1 import kraken
2
3 # Wczytanie obrazu
4 img = kraken.image.imread('tekst.jpg')
5
6 # Rozpoznawanie tekstu przy użyciu Kraken
7 text = kraken.recognize(img)
8
9 # Wypisanie rozpoznanego tekstu
10 print(text[0][1])
```

# EASYOCR

Szybkie i proste. Nie ma problemu z sieciami neuronowymi. Trochę słabiej działa na obrazach o niskiej jakości.

```
1 import easyocr
2
3 # Wczytanie obrazu
4 reader = easyocr.Reader(['en'])
5 results = reader.readtext('tekst.jpg')
6
7 # Wypisanie rozpoznanego tekstu
8 for res in results:
9     print(res[1])
```

## TESSERACT

W bibliotece Tesseract sam proces trenowania modeli OCR na bazie sieci neuronowych nie jest łatwy, ale możemy wykorzystać już gotowe modele, które zostały wytrenowane na dużych zbiorach danych. Domyślnie, Tesseract 4.0 wykorzystuje model LSTM (ang. Long Short-Term Memory) do rozpoznawania tekstu na obrazach. Aby uzyskać lepsze wyniki w optymalizacji OCR, należy wybrać odpowiedni język, a także dostosować odpowiednie opcje konfiguracyjne Tesseracta.

# TESSERACT

```
1 try:
2     from PIL import Image
3 except ImportError:
4     import Image
5 import pytesseract
6
7 # ścieżka do obrazka z tekstem
8 image_path = r'path\to\image.jpg'
9
10 # ustawienia konfiguracji OCR
11 custom_config = r'--oem 1 --psm 3'
12
13 # odczyt tekstu z obrazka przy użyciu Tesseract'a
14 text = pytesseract.image_to_string(Image.open(image_path), config=custom_config, lang='eng')
15
16 # wyświetlenie wyniku
17 print(text)
```

import os

1 obrazy/img.jpg



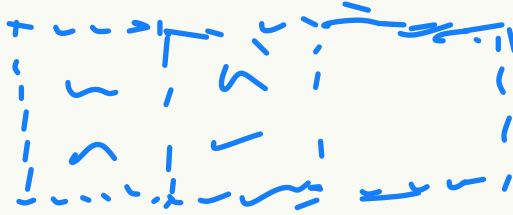
# KRAKEN

```
1 from kraken import pageseg
2 from kraken import binarization
3 from kraken import rnn
4 from kraken.lib import dataset
5 from kraken.lib.util import check_url
6
7 # pobieranie danych wyjściowych
8 datadir = 'https://github.com/mittagessen/kraken/raw/master/tests/data'
9 train_eq = check_url(f'{datadir}/00000014.eq.json')
10
11 # załadowanie zbiorów danych
12 with open(train_eq, 'r') as f:
13     data = dataset.build_dataset([train_eq])
14
15 # konfiguracja sieci neuronowej
16 model = rnn.build_model(train=data, network='lstm', layers=4, units=256, dropout=0.1)
17
18 # trenowanie sieci neuronowej
19 model = rnn.train_model(train=data, model=model, hidden=256, learn_rate=0.001, batch_size=3, max_iters=10)
20
21 # zapisanie wyuczonych wag sieci neuronowej
22 model.save_weights('trained_weights.h5')
```

*Handwritten notes:*

- A blue bracket and the word "importy" are written next to lines 1-5.
- A blue arrow points to the URL in line 8.
- A blue box highlights the parameters in line 19: `hidden=256, learn_rate=0.001, batch_size=3, max_iters=10`.

```
1 import easyocr
2 import json
3 import os
4
5 # Inicjalizacja czytnika z językiem angielskim
6 reader = easyocr.Reader(['en'])
7
8 # Ścieżka do obrazu
9 image_path = r'sciezka\do\obrazu.jpg'
10
11 # Odczytanie tekstu z obrazu przy użyciu OCR
12 results = reader.readtext(image_path)
13
14 # Inicjalizacja listy, która będzie przechowywać wpisy danych
15 dataset_entries = []
16
17 # Iteracja przez wyniki OCR i zapisanie tekstu oraz współrzędnych ramki
18 for (bbox, text, prob) in results:
19     data_entry = {
20         'tekst': text,
21         'prawdopodobienstwo': prob,
22         'ramka_obrysowujaca': bbox
23     }
24     dataset_entries.append(data_entry)
25
26 # Definicja ścieżki do zbioru danych
27 dataset_path = r'sciezka\do\zbioru_danych.json'
28
29 # Zapisanie wpisów danych do pliku JSON do późniejszego wykorzystania w uczeniu maszynowym
30 with open(dataset_path, 'w') as f:
31     json.dump(dataset_entries, f)
32
33 # Informacja dla użytkownika
34 print(f"Wyniki OCR zostały zapisane w {dataset_path}")
```



A hand-drawn diagram in blue ink showing a dashed rectangular box. Inside the box, there are several wavy lines representing text. This diagram illustrates the concept of a bounding box (bbox) used in OCR to identify and crop text from an image.